# QUALITY OF SERVICE · PART 1

## Quality of Service Models

**Best Effort** · No QoS policies are implemented

**Integrated Services (IntServ)**
Resource Reservation Protocol (RSVP) is used to reserve bandwidth per-flow across all nodes in a path

**Differentiated Services (DiffServ)**
Packets are individually classified and marked; policy decisions are made independently by each node in a path

## Layer 2 QoS Markings

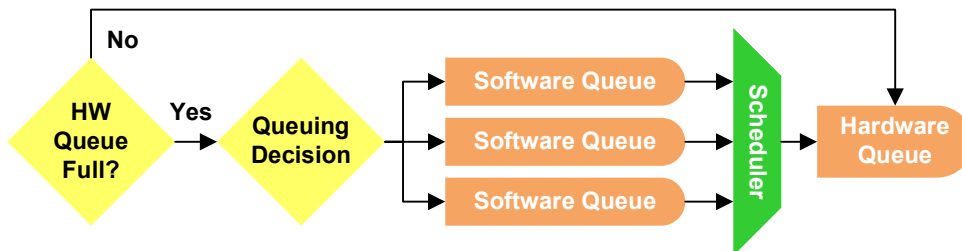| Medium | Name | Type |
|---|---|---|
| **Ethernet** | Class of Service (CoS) | 3-bit 802.1p field in 802.1Q header |
| **Frame Relay** | Discard Eligibility (DE) | 1-bit drop eligibility flag |
| **ATM** | Cell Loss Priority (CLP) | 1-bit drop eligibility flag |
| **MPLS** | Traffic Class (TC) | 3-bit field compatible with 802.1p |

## IP QoS Markings

**IP Precedence**
The first three bits of the IP TOS field; limited to 8 traffic classes

**Differentiated Services Code Point (DSCP)**
The first six bits of the IP TOS are evaluated to provide more granular classification; backward-compatible with IP Precedence

## QoS Flowchart



## Terminology

**Per-Hop Behavior (PHB)**
The individual QoS action performed at each independent DiffServ node

**Trust Boundary** · Beyond this, inbound QoS markings are not trusted

**Tail Drop** · Occurs when a packet is dropped because a queue is full

**Policing**
Imposes an artificial ceiling on the amount of bandwidth that may be consumed; traffic exceeding the policer rate is reclassified or dropped

**Shaping**
Similar to policing but buffers excess traffic for delayed transmission; makes more efficient use of bandwidth but introduces a delay

**TCP Synchronization**
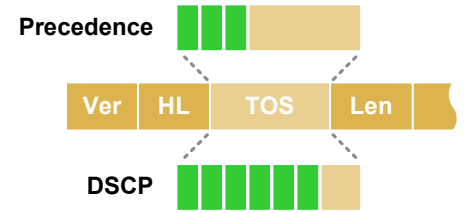Flows adjust TCP window sizes in synch, making inefficient use of a link

## DSCP Per-Hop Behaviors

**Class Selector (CS)** · Backward-compatible with IP Precedence values

**Assured Forwarding (AF)** · Four classes with variable drop preferences

**Expedited Forwarding (EF)** · Priority queuing for delay-sensitive traffic

## IP Type of Service (TOS)



## Precedence/DSCP

| | Binary | DSCP | Prec. |
|---|---|---|---|
| **56** | 111000 | Reserved | 7 |
| **48** | 110000 | Reserved | 6 |
| **46** | 101110 | EF | 5 |
| **32** | 100000 | CS4 | |
| **34** | 100010 | AF41 | 4 |
| **36** | 100100 | AF42 | |
| **38** | 100110 | AF43 | |
| **24** | 011000 | CS3 | |
| **26** | 011010 | AF31 | 3 |
| **28** | 011100 | AF32 | |
| **30** | 011110 | AF33 | |
| **16** | 010000 | CS2 | |
| **18** | 010010 | AF21 | 2 |
| **20** | 010100 | AF22 | |
| **22** | 010110 | AF23 | |
| **8** | 001000 | CS1 | |
| **10** | 001010 | AF11 | 1 |
| **12** | 001100 | AF12 | |
| **14** | 001110 | AF13 | |
| **0** | 000000 | BE | 0 |

## Congestion Avoidance

**Random Early Detection (RED)**
Packets are randomly dropped before a queue is full to prevent tail drop; mitigates TCP synchronization

**Weighted RED (WRED)**
RED with the added capability of recognizing prioritized traffic based on its marking

**Class-Based WRED (CBWRED)**
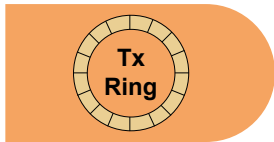WRED employed inside a class-based WFQ (CBWFQ) queue

# QUALITY OF SERVICE · PART 2

## Queuing Comparison

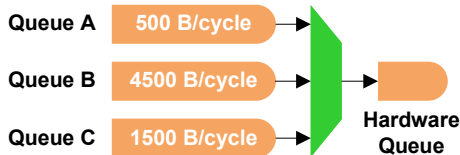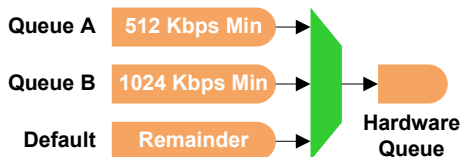| | FIFO | PQ | CQ | WFQ | CBWFQ | LLQ |
|---|---|---|---|---|---|---|
| **Default on Interfaces** | >2 Mbps | No | No | <=2 Mbps | No | No |
| **Number of Queues** | 1 | 4 | Configured | Dynamic | Configured | Configured |
| **Configurable Classes** | No | Yes | Yes | No | Yes | Yes |
| **Bandwidth Allocation** | Automatic | Automatic | Configured | Automatic | Configured | Configured |
| **Provides for Minimal Delay** | No | Yes | No | No | No | Yes |
| **Modern Implementation** | Yes | No | No | No | Yes | Yes |

## First In First Out (FIFO)



Hardware Queue

· Packets are transmitted in the order they are processed

· No prioritization is provided

· Default queuing method on high-speed (>2 Mbps) interfaces

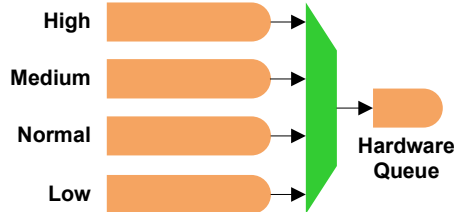· Configurable with the **tx-ring-limit** interface config command

## Custom Queuing (CQ)



· Rotates through queues using Weighted Round Robin (WRR)

· Processes a configurable number of bytes from each queue per turn

· Prevents queue starvation but does not provide for delay-sensitive traffic
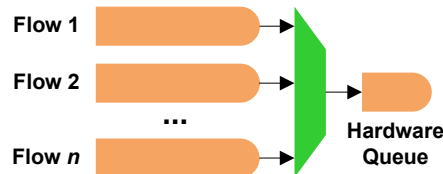
## Class-Based WFQ (CBWFQ)



· WFQ with administratively configured queues

· Each queue is allocated an amount/percentage of bandwidth

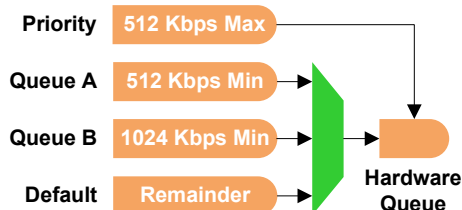· No support for delay-sensitive traffic

## Priority Queuing (PQ)



· Provides four static queues which cannot be reconfigured

· Higher-priority queues are always emptied before lower-priority queues

· Lower-priority queues are at risk of bandwidth starvation

## Weighted Fair Queuing (WFQ)



· Queues are dynamically created per flow to ensure fair processing

· Statistically drops packets from aggressive flows more often

· No support for delay-sensitive traffic

## Low Latency Queuing (LLQ)



· CBWFQ with the addition of a policed strict-priority queue

· Highly configurable while still supporting delay-sensitive traffic

## LLQ Config Example

```
                          Class Definitions
! Match packets by DSCP value
class-map match-all Voice
 match dscp ef
!
class-map match-all Call-Signaling
 match dscp cs3
!
class-map match-any Critical-Apps
 match dscp af21 af22
!
! Match packets by access list
class-map match-all Scavenger
 match access-group name Other
```

```
policy-map Foo              Policy Creation
 class Voice
  ! Priority queue policed to 33%
  priority percent 33
 class Call-Signaling
  ! Allocate 5% of bandwidth
  bandwidth percent 5
 class Critical-Apps
  bandwidth percent 20
  ! Extend queue size to 96 packets
  queue-limit 96
 class Scavenger
  ! Police to 64 kbps
  police cir 64000
   conform-action transmit
   exceed-action drop
 class class-default
  ! Enable WFQ
  fair-queue
  ! Enable WRED
  random-detect
```

```
interface Serial0          Policy Application
 ! Apply the policy in or out
 service-policy output Foo
```

## LLQ Config Example

```
show policy-map [interface]
```

```
Show interface
```

```
show queue <interface>
```

```
Show mls qos
```